CAN DEEPFAKE SPEECH BE RELIABLY DETECTED?

Anonymous authors

Paper under double-blind review

Abstract

Recent advances in text-to-speech (TTS) systems, particularly those with voice cloning capabilities, have made voice impersonation readily accessible, raising ethical and legal concerns due to potential misuse for malicious activities like misinformation campaigns and fraud. While synthetic speech detectors (SSDs) exist to combat this, they are vulnerable to "test domain shift", exhibiting decreased performance when audio is altered through transcoding, playback, or background noise. This vulnerability is further exacerbated by deliberate manipulation of synthetic speech aimed at deceiving detectors. This work presents the first systematic study of such active malicious attacks against state-of-the-art open-source SSDs. White-box attacks, black-box attacks, and their transferability are studied from both attack effectiveness and stealthiness, using both hardcoded metrics and human ratings. The results highlight the urgent need for more robust detection methods in the face of evolving adversarial threats.

022

004

006 007

008 009

010

011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

025 Recent years have witnessed a remarkable advance in text-to-speech (TTS) systems (Kreuk et al., 026 2022; Borsos et al., 2023; Leng et al., 2023; Saeki et al., 2023; Shen et al., 2023; Wang et al., 027 2023; ChatTTS, 2024; Chen et al., 2024; xTTS, 2024; ElevenLabs, 2024; Le et al., 2024; Lux et al., 028 2024). Some of these systems possess the zero-shot / few-shot voice cloning capability (Biadsy 029 et al., 2024; Cooper et al., 2020; Casanova et al., 2022; Ye et al., 2024; Le et al., 2024), that can mimic someone's voice using only a brief sample of that person's speech recording. The widespread 031 access of these systems, either through open-source projects or commercial APIs, has made voice impersonation easier than ever. This raises significant ethical and legal concerns as the capability can 033 be easily misused for misinformation campaign, fraud, copyright infringement, etc. As an instance, 034 a scammer utilized synthetic audio to mimic President Biden in unlawful robocalls during a New Hampshire primary election, resulting in a \$6 million penalty and felony accusations (Coldewey, 2024). Besides, numerous instances of synthetic audio misuses can be found online if searching with "deepfake audio", which underscores the urgent need to address this growing problem. 037

Synthetic speech detectors (SSDs) are deployed to mitigate the misuse of synthetic speech. While there are a number of advanced SSDs (Tak et al., 2021c;a; Jung et al., 2022), recent research (Müller et al., 2022; Xie et al., 2024) suggests they might struggle when facing "test domain shift". Con-040 cretely, a detector may exhibit decreased performance when presented with audio that has undergone 041 alterations, including transcoding, playback, background noise, or even just a shift in the TTS sys-042 tem used. However, current research efforts are directed towards these natural changes to audio, it is 043 important to recognize that deliberate manipulation of synthetic speech by an attacker with the intent 044 to deceive detectors can significantly increase the likelihood of success of the attacker. However, 045 there is a lack of systematic research on these malicious perturbations. 046

In this work, we conduct the first systematic study of active malicious attacks against the state-ofthe-art open-source SSDs. We investigate a range of attack scenarios, considering adversaries with
varying levels of access to the target SSD: those with full knowledge of the model (white-box), those
who can only interact with it and observe the results (black-box), and those who cannot even query
the model (agnostic); and evaluate the results using both hard-coded metrics and human ratings.

052 Our findings reveal that:

⁰⁵³

^{*}Equal Contribution

- Increased access to the detector makes it easier for attackers to create deepfakes that can evade detection without any noticeable loss in audio quality.
 Existing open-source SSD detectors are vulnerable when facing synthetic audio generated by TTS systems never seen during training.
 - VisQOL scores and human ratings show that the audio quality after attack is reasonable.
 - Alarmingly, even in the agnostic setting, attackers can still bypass state-of-the-art opensource SSDs with reasonable chance.

Overall, we need more robust SSDs to mitigate the growing threat of deepfake audio misuse.

063 064 065

066

056

059 060

061

062

2 PRELIMINARY

067 2.1 TTS TECHNIQUES

069 TTS systems, which convert written text into speech, have a long history of development and have 070 made remarkable progress in recent years. Early TTS systems primarily used a concatenative ap-071 proach (Khan & Chitode, 2016), where speech was synthesized by joining pre-recorded speech units 072 from a database. Despite the simplicity, they suffered from unnatural prosody and robotic-sounding 073 speech. To address the issue, researchers proposed statistical parametric speech synthesis (Zen et al., 2009). These systems used statistical models to learn the relationship between linguistic features 074 (e.g., phonemes, part-of-speech tags) and acoustic features (e.g., fundamental frequency, spectral 075 envelope) and enabled more natural-sounding speech generation with improved prosody. 076

Recently, deep learning has revolutionized the field of TTS. Neural network-based TTS systems have surpassed traditional methods in terms of speech naturalness and intelligibility. Several architectures have been explored, including sequence-to-sequence models (Wang et al., 2017; Ping et al., 2017), attention-based models (Ren et al., 2019; 2020), and generative adversarial networks (GANs) (Kumar et al., 2019). These models can directly learn the mapping from text to speech, enabling end-to-end training and eliminating the need for complex feature engineering. Some of the popular neural TTS systems include Kreuk et al. (2022); Borsos et al. (2023); Leng et al. (2023); Saeki et al. (2023); Shen et al. (2023); Wang et al. (2023); ChatTTS (2024); Chen et al. (2024); xTTS (2024); ElevenLabs (2024); Le et al. (2024); Lux et al. (2024).

085 086 087

2.2 Synthetic Speech Detection Techniques

In the past, detecting synthetic speech required carefully crafted features (Doddington et al., 2001; Alegre et al., 2013; Hanilçi et al., 2015; Patel & Patil, 2015; Sahidullah et al., 2015; Todisco et al., 2016). However, with the increase in available data and the development of larger models, simpler 091 features like waveforms or spectrograms are now sufficient for effective detection. RawNet2 (Tak 092 et al., 2021c) is a deep convolutional neural network for synthetic speech detection using merely raw waveforms. It builds upon the RawNet (Jung et al., 2020) architecture by incorporating residual 093 connections and dilated convolutions. RawNetGATST (Tak et al., 2021a) extends RawNet2 by 094 incorporating a graph attention network (Tak et al., 2021b) to identify key spectral or temporal 095 features for detection. Similarly, AASIST (Jung et al., 2022) refines the graph network architecture 096 further for improved synthetic speech detection. 097

098 099

100

3 BYPASS SYNTHETIC SPEECH DETECTION SYSTEMS

101 In this section, we aim to answer the following question:

Can deepfake audio be altered in ways nearly imperceptible to the human ear, but sufficient to bypass state-of-the-art detectors?

Unlike previous research that focused on natural perturbations (Müller et al., 2022; Xie et al., 2024),
we consider a malicious attacker who deliberately optimizes the perturbation to evade detection.
We examine this scenario under various levels of access to the detection systems, from having full knowledge (white-box), to partial knowledge (black-box), to no knowledge (agnostic).

108

161



Table 1: Baseline EERs of SSDs on the ASVSpoof2019-LA test split without attacks.

- White-box attacks can be highly effective and stealthy simultaneously.
 - 3



We first study white-box attack, where the adversary has full access to the model. We choose two white-box attacks: Projected Gradient Descent (PGD) (Mkadry et al., 2017) and I-FGSM (Kurakin et al., 2018). Algorithm 1 in Appendix shows details of PGD and I-FGSM.

Projected Gradient Descent: PGD crafts adversarial examples by iteratively taking small steps in the direction that maximizes the model's error, while projecting the perturbed example back within

ASVspoof WaveFake In-the-wild AASIST 0.970 ± 0.063 0.971 ± 0.046 0.985 ± 0.030 0.979 ± 0.037 AASIST-L 0.979 ± 0.045 0.975 ± 0.036 0.971 ± 0.077 RawNet2 1.000 ± 0.000 0.967 ± 0.063 RawGATST 0.997 ± 0.008 0.986 ± 0.030 0.997 ± 0.008

Table 2: Human ratings of speaker similarity between the original and PGD attacked audio.

a certain boundary around the original input to maintain a balance between attack success rate and stealthiness.

PGD has three major hyper-parameters: perturbation step size, ℓ_{∞} -norm constraint, and the number of iterations. We conduct hyper-parameter search and summarize the results in Figure 1, 2, and 3.

229 In Figure 1, we can tell that on WaveFake and In-the-wild, the attack success rate is almost always 230 100% while on ASVSpoof2019-LA test the attack success rate hovers between 60% and 100% 231 depending on the learning rate used. This reflects the fact that the detectors are more robust on 232 test data generated by the same TTS systems as the training data (*i.e.* in-domain data), but are still 233 vulnerable under white-box attacks with a few steps of hyper-parameter search. On the other hand, VisQOL scores keep decreasing as the perturbation step size grows. Usually, VisQOL score above 234 3.0 is considered reasonable quality. Thus, there exists a sweet spot of perturbation step size striking 235 balance between attack effectiveness and stealthiness. 236

In Figure 2, the observation of SSDs being more robust on ASVSpoof2019-LA test holds true. However, we observe that the VisQOL scores are pretty consistent despite the changing ℓ_{∞} -norm constraint, which says that audio quality is insensitive to ℓ_{∞} -norm constraint within a certain range.

Figure 3 shows that white-box attacks are efficient, reaching maximum attack success rates and
 stable VisQOL scores after just 50 iterations.

We also collect human ratings on whether the PGD-attacked audio with the best hyper-parameter combination sounds like the original synthetic audio, and the results are summarized in Table 2. We can see that most human raters think the two audio sound like the same person, underscoring the potential threat of using the attacked audio for impersonation.

246 247 248

249

250

251 252

253 254

255

256 257

258 259

260

261

262

263 264

216

217 218

219

220

222

Iterative Fast Gradient Sign Method: I-FGSM only differs from PGD in that it only uses the sign of the gradient to perturb the input audio. It shares the same set of hyper-parameters as PGD, for which the grid search results are summarized in Figure 4, 5, and 6 and human ratings are summarized in Table 3, and the findings are similar to PGD.

3.3 BLACK-BOX ATTACK

Takeaway:

- Existing open-source SSDs are still vulnerable to black-box attacks.
- Black-box attacks can be effective and stealthy simultaneously.

For black-box attack, we choose the Simple Black Box Attack (SimBA) (Guo et al., 2019). SimBA perturbs the input audio randomly and observes whether the prediction confidence score for "fake" class decreases or increases. If the confidence score decreases, SimBA will keep the perturbation. Otherwise, SimBA will try adding perturbation in the opposite direction and decide whether to

Table 3: Human ratings of speaker similarity between the original and I-FGSM attacked audio.

	ASVspoof	WaveFake	In-the-wild
AASIST	0.984 ± 0.020	0.960 ± 0.052	0.985 ± 0.024
AASIST-L	0.987 ± 0.022	0.986 ± 0.023	0.967 ± 0.054
RawNet2	0.980 ± 0.040	1.000 ± 0.000	0.991 ± 0.012
RawGATST	0.989 ± 0.024	0.858 ± 0.141	0.985 ± 0.024



- the details of SimBA.
- SimBA has three hyper-parameters: perturbation batch size, perturbation step size, and the number of queries. Perturbation batch size decides how many timesteps are perturbed in each query, while









	ASVspoof	WaveFake	In-the-wild
AASIST	0.984 ± 0.020	0.960 ± 0.052	0.985 ± 0.024
AASIST-L	0.987 ± 0.022	0.986 ± 0.023	0.967 ± 0.054
RawNet2	0.980 ± 0.040	1.000 ± 0.000	0.991 ± 0.012
RawGATST	0.989 ± 0.024	0.858 ± 0.141	0.985 ± 0.024

Table 4: Human ratings of speaker similarity between the original and simBA attacked audio.

Also in Figure 7, 8, and 9, we observe that ASSIST-L is the most robust model consistently, which is surprising because it's the smallest model within the 4 (See Jung et al. (2022) for the size of these models.). This observation aligns with the principle of Occam's razor, which suggests that simpler models often generalize better. A potential explanation could lie in the raggedness of the decision boundaries. Larger models, with their increased complexity, might create more intricate and potentially overfit decision boundaries. In contrast, ASSIST-L, being smaller, may form smoother decision boundaries, leading to better generalization and robustness against perturbations.

Human ratings of audio similarity is summarized in Table 4. Again the attacked audio sound highly similar to the original ones to human ears.

3.4 AGNOSTIC ATTACK: TRANSFERABILITY OF ABOVE ATTACKS

Takeaway:

- For both white-box attacks and black-box attacks, transferability depends on the target model's capability on the target audio.
- Black-box attacks are more transferrable on in-domain test data than out-of-domain data.
- Transferrability of different white-box attacks are alike.

The above attacks all assume different levels of access to the SSD model which might not be accessible in practice. As a result, we want to understand whether the above attacks are transferrable: Can a successfully attacked example on one model transfer to a different model? If this is true, then the adversary can craft a proxy model themselves, attack it, and expect it to bypass the real SSD as well.

The results are summarized in Figure 10. First, we find that on out-of-domain data, some SSDs are extremely vulnerable. For example, on WaveFake, RawNet2 is extremely vulnerable under all attacks; on In-the-wild, ASSIST and AASIST-L are more vulnerable than the other two models. Second, we find that on in-domain data, black-box attacks are much more transferrable than white-box attacks. This is because 1) black-box attacks tend to add larger perturbation than white-box attacks; 2) the SSDs' decision borders are alike for in-domain data. Thirdly, we also observe high similarity between the transferrability heatmap between PGD and I-FGSM, which might be due to different white-box attacks taking gradient paths in similar directions despite small differences.

4 CONCLUSION

This work presents the first systematic exploration of the robustness of state-of-the-art SSDs against
 adversarial attacks. Our findings reveal critical implications for SSD deployment and future re search.

Firstly, we demonstrate a clear correlation between system accessibility and vulnerability. Openaccess SSDs, and even those with oracle access, are highly susceptible to attacks. This underscores
the critical need to restrict public access to SSD models and internal workings. While complete
prevention of information leakage may be challenging, measures such as rate limiting can effectively
mitigate the threat of black-box attacks.



In conclusion, this study serves as a critical analysis of the current state of open-source SSD robust ness. By exposing key vulnerabilities and providing actionable recommendations, we aim to guide
 future research and development efforts towards building more secure and resilient SSD systems in
 an ever-evolving landscape of speech synthesis and spoofing technologies.

540 REFERENCES

549

- Federico Alegre, Ravichander Vipperla, Asmaa Amehraye, and Nicholas Evans. A new speaker verification spoofing countermeasure based on local binary patterns. In *INTERSPEECH 2013*, 14th Annual Conference of the International Speech Communication Association, Lyon: France (2013), pp. 5p, 2013.
- Fadi Biadsy, Youzheng Chen, Isaac Elias, Kyle Kastner, Gary Wang, Andrew Rosenberg, and Bhuvana Ramabhadran. Zero-shot cross-lingual voice transfer for tts. *arXiv preprint* arXiv:2409.13910, 2024.
- Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, et al. Audiolm: a language modeling approach to audio generation. *IEEE/ACM transactions on audio, speech, and language processing*, 31:2523–2533, 2023.
- Edresson Casanova, Julian Weber, Christopher D Shulby, Arnaldo Candido Junior, Eren Gölge, and
 Moacir A Ponti. Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for
 everyone. In *International Conference on Machine Learning*, pp. 2709–2720. PMLR, 2022.
- ChatTTS. Chattts text-to-speech for conversational scenarios. https://chattts.com/, 2024. Accessed: 07/16/2024.
- Sanyuan Chen, Shujie Liu, Long Zhou, Yanqing Liu, Xu Tan, Jinyu Li, Sheng Zhao, Yao Qian, and
 Furu Wei. Vall-e 2: Neural codec language models are human parity zero-shot text to speech
 synthesizers. *arXiv preprint arXiv:2406.05370*, 2024.
- 563
 564
 564
 565
 566
 566
 567
 568
 568
 569
 569
 560
 560
 560
 561
 562
 562
 563
 564
 564
 565
 566
 566
 566
 566
 566
 567
 568
 568
 568
 569
 560
 560
 560
 560
 560
 561
 562
 562
 564
 565
 566
 566
 566
 566
 566
 566
 567
 568
 568
 568
 568
 568
 569
 569
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
 560
- Erica Cooper, Cheng-I Lai, Yusuke Yasuda, Fuming Fang, Xin Wang, Nanxin Chen, and Junichi
 Yamagishi. Zero-shot multi-speaker text-to-speech with state-of-the-art neural speaker embeddings. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6184–6188. IEEE, 2020.
- George R Doddington et al. Speaker recognition based on idiolectal differences between speakers. In *Interspeech*, pp. 2521–2524, 2001.
- 574 ElevenLabs. High quality, human-like ai voice generator. https://elevenlabs.io/ 575 text-to-speech, 2024. Accessed: 07/16/2024.
- Joel Frank and Lea Schönherr. Wavefake: A data set to facilitate audio deepfake detection. *arXiv* preprint arXiv:2111.02813, 2021.
- 579 Chuan Guo, Jacob Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Weinberger. Simple
 580 black-box adversarial attacks. In *International conference on machine learning*, pp. 2484–2493.
 581 PMLR, 2019.
- Cemal Hanilçi, Tomi Kinnunen, Md Sahidullah, and Aleksandr Sizov. Classifiers for synthetic speech detection: A comparison. 2015.
- Andrew Hines, Jan Skoglund, Anil C Kokaram, and Naomi Harte. Visqol: an objective speech quality model. *EURASIP Journal on Audio, Speech, and Music Processing*, 2015:1–18, 2015.
- Jee-weon Jung, Seung-bin Kim, Hye-jin Shim, Ju-ho Kim, and Ha-Jin Yu. Improved rawnet with feature map scaling for text-independent speaker verification using raw waveforms. *arXiv preprint arXiv:2004.00526*, 2020.
- Jee-weon Jung, Hee-Soo Heo, Hemlata Tak, Hye-jin Shim, Joon Son Chung, Bong-Jin Lee, Ha-Jin
 Yu, and Nicholas Evans. Aasist: Audio anti-spoofing using integrated spectro-temporal graph
 attention networks. In *ICASSP 2022-2022 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pp. 6367–6371. IEEE, 2022.

594 Rubeena A Khan and Janardan Shrawan Chitode. Concatenative speech synthesis: A review. International Journal of Computer Applications, 136(3):1–6, 2016. 596 Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi 597 Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. arXiv 598 preprint arXiv:2209.15352, 2022. 600 Kundan Kumar, Rithesh Kumar, Thibault De Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, 601 Alexandre De Brebisson, Yoshua Bengio, and Aaron C Courville. Melgan: Generative adversarial 602 networks for conditional waveform synthesis. Advances in neural information processing systems, 603 32, 2019. 604 Alexey Kurakin, Ian J Goodfellow, and Samy Bengio. Adversarial examples in the physical world. 605 In Artificial intelligence safety and security, pp. 99–112. 2018. 606 607 Matthew Le, Apoorv Vyas, Bowen Shi, Brian Karrer, Leda Sari, Rashel Moritz, Mary Williamson, 608 Vimal Manohar, Yossi Adi, Jay Mahadeokar, et al. Voicebox: Text-guided multilingual universal 609 speech generation at scale. Advances in neural information processing systems, 36, 2024. 610 611 Yichong Leng, Zhifang Guo, Kai Shen, Xu Tan, Zeqian Ju, Yanqing Liu, Yufei Liu, Dongchao Yang, Leying Zhang, Kaitao Song, et al. Prompttts 2: Describing and generating voices with text 612 prompt. arXiv preprint arXiv:2309.02285, 2023. 613 614 Hongbin Liu, Moyang Guo, Zhengyuan Jiang, Lun Wang, and Neil Zhenqiang Gong. Audiomark-615 bench: Benchmarking robustness of audio watermarking. arXiv preprint arXiv:2406.06979, 2024. 616 617 Florian Lux, Sarina Meyer, Lyonel Behringer, Frank Zalkow, Phat Do, Matt Coler, Emanuël A. P. Habets, and Ngoc Thang Vu. Meta Learning Text-to-Speech Synthesis in over 7000 Languages. 618 In Interspeech. ISCA, 2024. 619 620 Aleksander Mkadry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 621 Towards deep learning models resistant to adversarial attacks. *stat*, 1050(9), 2017. 622 623 Nicolas M Müller, Pavel Czempin, Franziska Dieckmann, Adam Froghyar, and Konstantin 624 Böttinger. Does audio deepfake detection generalize? arXiv preprint arXiv:2203.16263, 2022. 625 Tanvina B Patel and Hemant A Patil. Combining evidences from mel cepstral, cochlear filter cepstral 626 and instantaneous frequency features for detection of natural vs. spoofed speech. In Interspeech, 627 pp. 2062–2066, 2015. 628 629 Wei Ping, Kainan Peng, Andrew Gibiansky, Sercan O Arik, Ajay Kannan, Sharan Narang, Jonathan 630 Raiman, and John Miller. Deep voice 3: Scaling text-to-speech with convolutional sequence 631 learning. arXiv preprint arXiv:1710.07654, 2017. 632 Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech: 633 Fast, robust and controllable text to speech. Advances in neural information processing systems, 634 32, 2019. 635 636 Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech 2: 637 Fast and high-quality end-to-end text to speech. arXiv preprint arXiv:2006.04558, 2020. 638 Takaaki Saeki, Heiga Zen, Zhehuai Chen, Nobuyuki Morioka, Gary Wang, Yu Zhang, Ankur Bapna, 639 Andrew Rosenberg, and Bhuvana Ramabhadran. Virtuoso: Massive multilingual speech-text joint 640 semi-supervised learning for text-to-speech. In ICASSP 2023-2023 IEEE International Confer-641 ence on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5. IEEE, 2023. 642 643 Md Sahidullah, Tomi Kinnunen, and Cemal Hanilçi. A comparison of features for synthetic speech 644 detection. 2015. 645 Kai Shen, Zeqian Ju, Xu Tan, Yanqing Liu, Yichong Leng, Lei He, Tao Qin, Sheng Zhao, and Jiang 646 Bian. Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing 647 synthesizers. arXiv preprint arXiv:2304.09116, 2023.

- Hemlata Tak, Jee-weon Jung, Jose Patino, Madhu Kamble, Massimiliano Todisco, and Nicholas
 Evans. End-to-end spectro-temporal graph attention networks for speaker verification anti spoofing and speech deepfake detection. *arXiv preprint arXiv:2107.12710*, 2021a.
- Hemlata Tak, Jee-weon Jung, Jose Patino, Massimiliano Todisco, and Nicholas Evans. Graph atten tion networks for anti-spoofing. *arXiv preprint arXiv:2104.03654*, 2021b.
- Hemlata Tak, Jose Patino, Massimiliano Todisco, Andreas Nautsch, Nicholas Evans, and Anthony
 Larcher. End-to-end anti-spoofing with rawnet2. In *ICASSP 2021-2021 IEEE International Con- ference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6369–6373. IEEE, 2021c.
- Massimiliano Todisco, Héctor Delgado, and Nicholas WD Evans. A new feature for automatic speaker verification anti-spoofing: Constant q cepstral coefficients. In *Odyssey*, volume 2016, pp. 283–290, 2016.
- Massimiliano Todisco, Xin Wang, Ville Vestman, Md Sahidullah, Héctor Delgado, Andreas
 Nautsch, Junichi Yamagishi, Nicholas Evans, Tomi Kinnunen, and Kong Aik Lee. Asvspoof
 2019: Future horizons in spoofed and fake audio detection. *arXiv preprint arXiv:1904.05441*, 2019.
- Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. Neural codec language models are zero-shot text to speech synthesizers. *arXiv preprint arXiv:2301.02111*, 2023.
- Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly,
 Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. Tacotron: Towards end-to-end
 speech synthesis. *arXiv preprint arXiv:1703.10135*, 2017.
- Yuankun Xie, Yi Lu, Ruibo Fu, Zhengqi Wen, Zhiyong Wang, Jianhua Tao, Xin Qi, Xiaopeng Wang,
 Yukun Liu, Haonan Cheng, et al. The codecfake dataset and countermeasures for the universally
 detection of deepfake audio. *arXiv preprint arXiv:2405.04880*, 2024.
- artts. High quality, human-like ai voice generator. https://github.com/coqui-ai/TTS, 2024. Accessed: 07/16/2024.
- ⁶⁷⁸ Zhen Ye, Zeqian Ju, Haohe Liu, Xu Tan, Jianyi Chen, Yiwen Lu, Peiwen Sun, Jiahao Pan, Weizhen
 ⁶⁷⁹ Bian, Shulin He, et al. Flashspeech: Efficient zero-shot speech synthesis. *arXiv preprint*⁶⁸⁰ *arXiv:2404.14700*, 2024.
- Heiga Zen, Keiichi Tokuda, and Alan W Black. Statistical parametric speech synthesis. *speech communication*, 51(11):1039–1064, 2009.

$\begin{aligned} \varepsilon_{\text{reactions } f, \ f \in a \ \text{durb class } \mathcal{K}, \ \text{attack } A \in \{ \text{ PGD}, \ \text{FPGSM} \}, \ \varepsilon_{\infty}\text{-from constraint } \mathcal{E} \\ \text{Output: Adversarial perturbation } \delta \in \mathbb{R}^T \\ 1: \ \delta \leftarrow 0 \\ \text{s} \text{if } A = \text{*PGD' then} \\ 4: \delta \leftarrow \delta - \alpha \cdot \nabla_{\delta} l_{CE}(f(s + \delta), \mathcal{R}) \\ \text{s} \text{end if} \\ 6: \text{if } A = \text{*I-FGSM' then} \\ \text{?} \delta \leftarrow \delta - \alpha \cdot \text{sign}(\nabla_{\delta} l_{CE}(f(s + \delta), \mathcal{R})) \\ \text{s} \text{end if} \\ 0: \text{if } f(s + \delta) = = \mathcal{R} \ \text{then} \\ \text{begin{subarray}{l} b \ \text{constraint } b \ \text{oprojection} \\ \text{b } \delta \leftarrow \text{clip}(\delta, -\epsilon, \epsilon) \\ \text{clip}(\delta, -\epsilon, \epsilon) \\ \text{l} \text{constraint } b \ \text{constraint } \epsilon \\ \text{ottack success} \\ 1: \text{Break} \\ 12: \text{end if} \\ 13: \ \text{end for} \\ 14: \ \text{returb } \delta \\ \hline \\$;_	: Waveform audio $s \in \mathbb{R}^{1}$, SSD	model f, perturbation step size α , maximum number of $k \in [(PGD)^2 (I EGSM)]^{\ell}$ norm constraint.
$ s \ def{alpha} = \delta \ def{alpha} s \ def{alpha} s$	11 Outn	erations 1, real audio class \mathcal{K} , attac	$\kappa A \in \{ \text{ PGD}, \text{ I-FGSM} \}, \ell_{\infty}$ -norm constraint ϵ
$\begin{array}{cccccccccccccccccccccccccccccccccccc$	0 uιμ 1 · δ		⊳ initializatior
2. if A == PGD then 4: $\delta \leftarrow \delta - \alpha \cdot \nabla_{\delta} l_{CE}(f(s + \delta), \mathcal{R})$ > gradient descent with cross entropy loss l_{CE} 5: end if 7: $\delta \leftarrow \delta - \alpha \cdot \operatorname{sign}(\nabla_{\delta} l_{CE}(f(s + \delta), \mathcal{R}))$ > FGSM with cross entropy loss l_{CE} 8: end if > projection 9: $\delta \leftarrow \operatorname{clip}(\delta, -\epsilon, \epsilon)$ > projection 10: if $f(s + \delta) == \mathcal{R}$ then > attack success 11: Break > attack 12: end if = attack 13: end for = attack 14: return δ = attack Augorithm 2 Black-box Attack: SimBA Intent is a success 11: Break = attack 12: end if = attack 13: end for = attack 14: return δ = attack Output: Waveform audio $s \in \mathbb{R}^T$, SD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ > initializing highest probability to predict real audio class 2: while t < T do > initialize highest probability to predict real audio class 3: while t < T do > attack success <	1. 0 2. f	- 0 or $t \leftarrow 1$ to T do	
$ \begin{array}{llllllllllllllllllllllllllllllllllll$	2. I 3.	if $A == PGD'$ then	
$ \begin{array}{llllllllllllllllllllllllllllllllllll$	3. 4·	$\delta \leftarrow \delta - \alpha \cdot \nabla sl_{CE}(f(s+\delta))$	\mathcal{R}) \triangleright gradient descent with cross entropy loss l_{CE}
6: if A == '1-FGSM' then 7: δ ← δ − α · sign(∇δl _{CE} (f(s + δ), R))) ▷ FGSM with cross entropy loss l _{CE} 8: end if ▷ projection 9: δ ← clip(δ, -ε, ε) ▷ projection 10: if f(s + δ) == R then ▷ attack success 11: Break ▷ attack success 12: end if □ 13: end for □ 14: return δ □ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f, perturbation step size α, perturbation batch size q maximum number of queries Q, real audio class R, ℓ _∞ -norm constraint ϵ Output: Adversarial perturbation δ ∈ \mathbb{R}^T 1: δ ← 0, t ← 0 ▷ initializing highest probability to predict real audio class 2: p ← f(s, R), t ← t + 1 ▷ initializing highest probability to predict real audio class 3: while t < T do	5:	end if	(,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,
7: $\delta \leftarrow \delta - \alpha \cdot \operatorname{sign}(\nabla_{\delta} l_{CE}(f(s + \delta), \mathcal{R}))$ FGSM with cross entropy loss l_{CE} 8: end if 9: $\delta \leftarrow \operatorname{clip}(\delta, -\epsilon, \epsilon)$ 10: if $f(s + \delta) == \mathcal{R}$ then 11: Break 12: end if 13: end for 14: return δ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ 3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ 11: $\delta = 0, t \leftarrow 0$ 11: $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then 13: $Continue$ 14: else 15: $t \leftarrow t + 1$ 16: $t \leftarrow t + 1$ 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 19: end if 20: end if 21: end while 22: return δ	6:	if $A ==$ 'I-FGSM' then	
8: end if 9: $\delta \leftarrow \operatorname{clip}(\delta, -\epsilon, \epsilon)$ 10: if $f(s + \delta) == \mathcal{R}$ then 11: Break 12: end if 13: end for 14: return δ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class $\mathcal{R}, \ell_{\infty}$ -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ 1: $\delta \leftarrow 0, t \leftarrow 1$ 1: $\delta \leftarrow 0, t \leftarrow 0$ 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ 1: $\delta \leftarrow 0, t \leftarrow 0$ 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ 1: $\delta = \mathcal{R}$ then 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ 12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then 13: Continue 14: else 15: $t \leftarrow t + 1$ 15: $t \leftarrow t + 1$ 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 18: Continue 19: end if 20: end if 21: end while 22: return δ	7:	$\delta \leftarrow \delta - \alpha \cdot \operatorname{sign}(\nabla_{\delta} l_{CE}) (f(s))$	$(+\delta), \mathcal{R})) $ \triangleright FGSM with cross entropy loss l_{CE}
9: $\delta \leftarrow \operatorname{clip}(\delta, -\epsilon, \epsilon)$ 10: if $f(s + \delta) == \mathbb{R}$ then 11: Break 12: end if 13: end for 14: return δ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class $\mathcal{R}, \ell_{\infty}$ -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ 3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ 1: $\phi - f(s + \delta + r, \mathcal{R}) > p$ then 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then 13: Continue 14: else 15: $t \leftarrow t + 1$ 16: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 19: end if 20: end if 21: end while 22: return δ	8:	end if	,, ,, ,, ,, ,, ,, ,, ,, ,, ,, ,, ,, ,,
10: if $f(s + \delta) == \mathcal{R}$ then ▷ attack success 11: Break ▷ 12: end if	9:	$\delta \leftarrow \operatorname{clip}(\delta, -\epsilon, \epsilon)$	▷ projectior
11: Break 12: end if 13: end for 14: return δ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ \triangleright initialization 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ \triangleright initializing highest probability to predict real audio class 3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then \triangleright attack success 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 18: Continue 19: end if 20: end if 21: end while 22: return δ	10:	if $f(s+\delta) == \mathcal{R}$ then	⊳ attack success
12: end fr 13: end for 14: return δ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ \triangleright initialization 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ \triangleright initializing highest probability to predict real audio class 3: while $t < T$ do \triangleright initializing highest probability to predict real audio class 3: while $t < T$ do \triangleright attack success 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	11:	Break	
13: end for 14: return δ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ ▷ initialization 1: $\delta \leftarrow 0, t \leftarrow 0$ ▷ initializing highest probability to predict real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ ▷ initializing highest probability to predict real audio class \mathcal{R} if $f(s + \delta) = \mathcal{R}$ then 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ ▷ initializing highest probability to predict real audio class \mathcal{R} 3: while $t < T$ do ▷ attack success 5: Break ▷ attack success 6: end if \mathcal{R} 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ ▷ attack success 8: Randomly choose q dimensions from r without replacement ○ one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then ▷ one more query for f below 12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then ▷ one more query for f below 13: Continue ▷ one more query for f below 14: else ▷ one more query for f below 15: $t \leftarrow t + 1$ ▷ one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) >$	12:	end if	
14: return δ Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ \triangleright initialization 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ \triangleright initializing highest probability to predict real audio class 3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then \triangleright attack success 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 18: Continue 19: end if 20: end if 21: end while 22: return δ	13: e	nd for	
Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ > initialization 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ > initializing highest probability to predict real audio class 3: while $t < T$ do > attack success 4: if $f(s + \delta) == \mathcal{R}$ then > attack success 5: Break > attack success 6: end if > $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement > one more query for f below 1: $f(s + \delta + r, \mathcal{R}) > p$ then > one more query for f below 1: if $f(s + \delta + r, \mathcal{R}) > p$ then > update highest probability to predict real audio class 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ > update highest probability to predict real audio class 13: Continue > one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then > update highest probability to predict real audio class 18: Continue > update highest probability to predict real audio class 19: end if 21: end while 22: return δ	14: r	eturn δ	
Algorithm 2 Black-box Attack: SimBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ \triangleright initialization $\delta \in \mathbb{R}^T$ $2: \phi \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ \triangleright initializing highest probability to predict real audio class $2: p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ \triangleright initializing highest probability to predict real audio class $3:$ while $t < T$ do \triangleright initializing highest probability to predict real audio class $3:$ while $t < T$ do \triangleright initializing highest probability to predict real audio class $3:$ while $t < T$ do \triangleright initializing highest probability to predict real audio class $3:$ while $t < T$ do \triangleright attack success $5:$ Break \triangleright attack success $6:$ end if $?$ $r \in \mathbb{R}^T$ and $r \leftarrow 0$ $8:$ Randomly choose q dimensions from r without replacement $?$ $?$ one more query for f below $10:$ $t \leftarrow t + 1$ \triangleright one more query for f below $?$ $?$ one more query for f below $11:$ $if f(s + \delta - r, \mathcal{R}) > p$ then $>$ one more query for f below $?$ $?$ one more query for f below $?$ $?$ $>$ o			
Argorithm 2 Black-box Attack: SIMBA Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Dutput: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ \triangleright initializing highest probability to predict real audio class 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ \triangleright initializing highest probability to predict real audio class 3: while $t < T$ do \triangleright attack success 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	A 1	the 2 Diast has Attack Sim DA	
Input: Waveform audio $s \in \mathbb{R}^T$, SSD model f , perturbation step size α , perturbation batch size q maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ \triangleright initialization 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ \triangleright initializing highest probability to predict real audio class 3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then \triangleright attack success 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 20: end if 21: end while 22: return δ	Algo	runm 2 Black-box Attack: SIMBA	
maximum number of queries Q , real audio class \mathcal{R} , ℓ_{∞} -norm constraint ϵ Output: Adversarial perturbation $\delta \in \mathbb{R}^T$ 1: $\delta \leftarrow 0, t \leftarrow 0$ > initialization2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ > initializing highest probability to predict real audio class3: while $t < T$ do> initializing highest probability to predict real audio class4: if $f(s + \delta) == \mathcal{R}$ then> attack success5: Break> end if6: end if> $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ > one more query for f below11: if $f(s + \delta + r, \mathcal{R}) > p$ then> one more query for f below12: $p \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then> one more query for f below13: ContinueIf $f(s + \delta - r, \mathcal{R}) > p$ then14: else> one more query for f below15: $t \leftarrow t + 1$ > update highest probability to predict real audio class16: if $f(s + \delta - r, \mathcal{R}) > p$ then> update highest probability to predict real audio class17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then> update highest probability to predict real audio class18: ContinueIf19: end if20: end if20: end if21: end while22: return δ	Inpu	: Waveform audio $s \in \mathbb{R}^T$, SSD m	odel f, perturbation step size α , perturbation batch size q
Output: Adversarial perturbation $\delta \in \mathbb{R}^{T}$ 1: $\delta \leftarrow 0, t \leftarrow 0$ > initialization 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ > initializing highest probability to predict real audio class 3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then > attack success 5: Break 6: end if 7: $r \in \mathbb{R}^{T}$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ > one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ > p then 13: Continue 14: else 15: $t \leftarrow t + 1$ > one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ > p then 18: Continue 19: end if 20: end if 21: end while 22: return δ	'n	naximum number of queries Q , real	audio class $\mathcal{R}, \ell_{\infty}$ -norm constraint ϵ
1: $\delta \leftarrow 0, t \leftarrow 0$ > initialization 2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ > initializing highest probability to predict real audio class 3: while $t < T$ do > attack success 4: if $f(s + \delta) == \mathcal{R}$ then > attack success 5: Break > attack success 6: end if ? $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement > one more query for f below 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r > one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then > one more query for f below 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ > update highest probability to predict real audio class 13: Continue > one more query for f below 14: else > one more query for f below 15: $t \leftarrow t + 1$ > one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then > update highest probability to predict real audio class 18: Continue > update highest probability to predict real audio class 19: end if 20: end if 21: end while 22: return δ	Outp	ut: Adversarial perturbation $\delta \in \mathbb{R}^{2}$	Γ
2: $p \leftarrow f(s, \mathcal{R}), t \leftarrow t+1$ initializing highest probability to predict real audio class 3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t+1$ 10: $t \leftarrow t+1$ 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ 13: Continue 14: else 15: $t \leftarrow t+1$ 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ 18: Continue 19: end if 20: end if 21: end while 22: return δ	1: δ	$\leftarrow 0, t \leftarrow 0$	
3: while $t < T$ do 4: if $f(s + \delta) == \mathcal{R}$ then \triangleright attack success 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 18: Continue 19: end if 20: end if 21: end while 22: return δ	-		\triangleright initialization
4: if $f(s + \delta) == \mathcal{R}$ then 5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ 10: $t \leftarrow t + 1$ 10: $t \leftarrow t + 1$ 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ 13: Continue 14: else 15: $t \leftarrow t + 1$ 16: $t \leftarrow t + 1$ 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p	$\leftarrow f(s,\mathcal{R}), t \leftarrow t+1 \qquad \qquad \texttt{I}$	▷ initialization > initializing highest probability to predict real audio class
5: Break 6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: v	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ this $t < T$ do	initialization initializing highest probability to predict real audio class
6: end if 7: $r \in \mathbb{R}^T$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R}) > p$ then 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then	 initialization initializing highest probability to predict real audio class attack success
7: $r \in \mathbb{R}^{r}$ and $r \leftarrow 0$ 8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t+1$ $\qquad \triangleright$ one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ $\qquad \triangleright$ update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t+1$ $\qquad \triangleright$ one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ $\qquad \triangleright$ update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5:	$ \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1 $ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break	 ▷ initialization > initializing highest probability to predict real audio class ▷ attack success
8: Randomly choose q dimensions from r without replacement 9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t + 1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if \mathcal{R}^T	 ▷ initialization > initializing highest probability to predict real audio class ▷ attack success
9: Randomly add α or $-\alpha$ to the chosen q dimensions in r 10: $t \leftarrow t+1$ \triangleright one more query for f below 11: if $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t+1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6: 7:	$ \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1 $ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$	 initializing highest probability to predict real audio class attack success
10: $t \leftarrow t + 1$ \triangleright one more query for f below11:if $f(s + \delta + r, \mathcal{R}) > p$ then \triangleright update highest probability to predict real audio class12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class13: Continue 14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below16: $if f(s + \delta - r, \mathcal{R}) > p$ then17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ 18: Continue 19:end if20:end if21:end while22:return δ	2: p 3: w 4: 5: 6: 7: 8:	$ \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1 $ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f	> initializing highest probability to predict real audio class > attack success rom r without replacement
11: If $f(s + \delta + r, \mathcal{R}) > p$ then 12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ 13: Continue 14: else 15: $t \leftarrow t + 1$ 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the choice	> initializing highest probability to predict real audio class > attack success rom r without replacement osen q dimensions in r
12: $p \leftarrow f(s + \delta + r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 13: Continue 14: else 15: $t \leftarrow t + 1$ 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9: 10:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ the t < T do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the choice $t \leftarrow t + 1$ if $f(\alpha + \delta) = \mathcal{R}$	> initializing highest probability to predict real audio class > attack success from r without replacement osen q dimensions in r > one more query for f below
14: else 15: $t \leftarrow t + 1$ \triangleright one more query for f below 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9: 10: 11:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the cho- $t \leftarrow t + 1$ if $f(s + \delta + r, \mathcal{R}) > p$ then $m \leftarrow f(s + \delta + r, \mathcal{R}) > p$ then	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class > attack success from r without replacement osen q dimensions in r > one more query for f below
15: $t \leftarrow t + 1$ 16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: v 4: 5: 6: 7: 8: 9: 10: 11: 12: 13:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ therefore the set of the se	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class from r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class
16: if $f(s + \delta - r, \mathcal{R}) > p$ then 17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: v 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14:	$ \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1 $ the conduct of the second state of th	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class rom r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class
17: $p \leftarrow f(s + \delta - r, \mathcal{R})$ \triangleright update highest probability to predict real audio class 18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: v 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the cho- $t \leftarrow t + 1$ if $f(s + \delta + r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta + r, \mathcal{R})$ Continue else $t \leftarrow t + 1$	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class from r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class > one more query for f below
18: Continue 19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the cho- $t \leftarrow t + 1$ if $f(s + \delta + r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta + r, \mathcal{R})$ Continue else $t \leftarrow t + 1$ if $f(s + \delta - r, \mathcal{R}) > n$ then	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class from r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class > one more query for f below
19: end if 20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16: 17:	$\leftarrow f(s, \mathcal{R}), t \leftarrow t + 1$ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the choice $t \leftarrow t + 1$ if $f(s + \delta + r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta + r, \mathcal{R})$ Continue else $t \leftarrow t + 1$ if $f(s + \delta - r, \mathcal{R}) > p$ then $n \leftarrow f(s + \delta - r, \mathcal{R})$	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class rom r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class > update highest probability to predict real audio class
20: end if 21: end while 22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16: 17: 18:	$ \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1 $ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the choose $t \leftarrow t + 1$ if $f(s + \delta + r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta + r, \mathcal{R})$ Continue else $t \leftarrow t + 1$ if $f(s + \delta - r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta - r, \mathcal{R})$ Continue	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class from r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class > one more query for f below > update highest probability to predict real audio class
21: end while 22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16: 17: 18: 19:	$ \begin{array}{l} \leftarrow f(s,\mathcal{R}), t\leftarrow t+1 & \text{tr}\\ \text{thile } t < T \text{ do} & \text{if } f(s+\delta) == \mathcal{R} \text{ then} \\ & \text{Break} & \text{end if} \\ r \in \mathbb{R}^T \text{ and } r \leftarrow 0 & \text{Randomly choose } q \text{ dimensions ff} \\ \text{Randomly add } \alpha \text{ or } -\alpha \text{ to the choing } t\leftarrow t+1 & \text{if } f(s+\delta+r,\mathcal{R}) > p \text{ then} \\ p\leftarrow f(s+\delta+r,\mathcal{R}) > p \text{ then} & p\leftarrow f(s+\delta+r,\mathcal{R}) & \text{Continue} \\ \text{else} & t\leftarrow t+1 & \text{if } f(s+\delta-r,\mathcal{R}) > p \text{ then} & p\leftarrow f(s+\delta-r,\mathcal{R}) & \text{Continue} \\ & \text{end if} & \text{continue} & \text{end if} & \text{continue} \\ \end{array} $	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class from r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class > one more query for f below > update highest probability to predict real audio class
22: return δ	2: p 3: w 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16: 17: 18: 19: 20:	$ \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1 $ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the choose $t \leftarrow t + 1$ if $f(s + \delta + r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta + r, \mathcal{R})$ Continue else $t \leftarrow t + 1$ if $f(s + \delta - r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta - r, \mathcal{R})$ Continue end if end if	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class from r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class > one more query for f below > update highest probability to predict real audio class
	2: p 3: w 4: 5: 6: 7: 8: 9: 10: 11: 12: 13: 14: 15: 16: 17: 18: 19: 20: e	$ \leftarrow f(s, \mathcal{R}), t \leftarrow t + 1 $ thile $t < T$ do if $f(s + \delta) == \mathcal{R}$ then Break end if $r \in \mathbb{R}^T$ and $r \leftarrow 0$ Randomly choose q dimensions f Randomly add α or $-\alpha$ to the choose $t \leftarrow t + 1$ if $f(s + \delta + r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta + r, \mathcal{R})$ Continue else $t \leftarrow t + 1$ if $f(s + \delta - r, \mathcal{R}) > p$ then $p \leftarrow f(s + \delta - r, \mathcal{R})$ Continue end if end if end if nd while	> initializing highest probability to predict real audio class > initializing highest probability to predict real audio class from r without replacement osen q dimensions in r > one more query for f below > update highest probability to predict real audio class > one more query for f below > update highest probability to predict real audio class